

Solving L_0 -Problems via Mixed-Integer Optimization

Théo Guyard

Inria, Centre de l'Université de Rennes, France

Insa Rennes, IRMAR CNRS UMR 6625, France

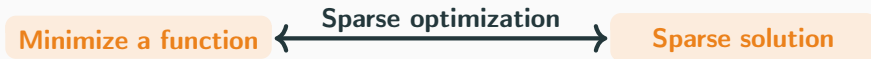
LS2N seminar

7th of March, 2024

Nantes, France

Sparse Optimization

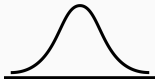
Two goals, one problem



Signal processing



Machine learning

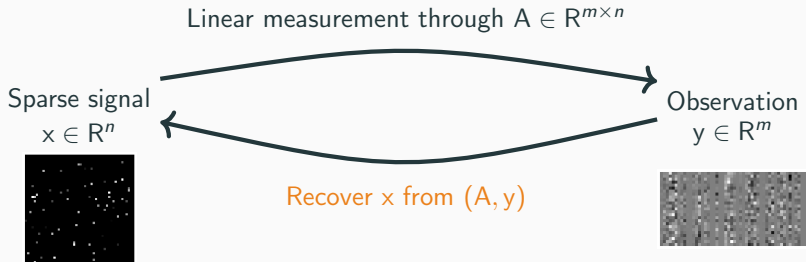


High-dim. statistics



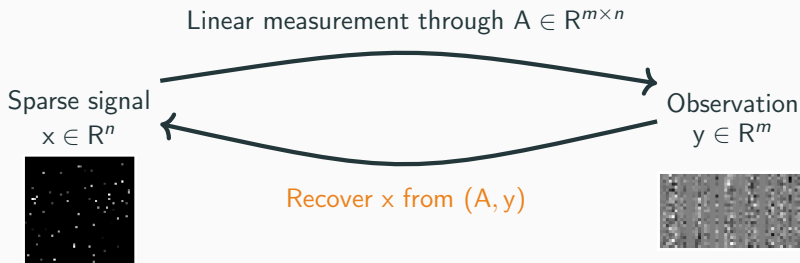
And many others

Compressive sensing



Find x such that $y \simeq Ax$

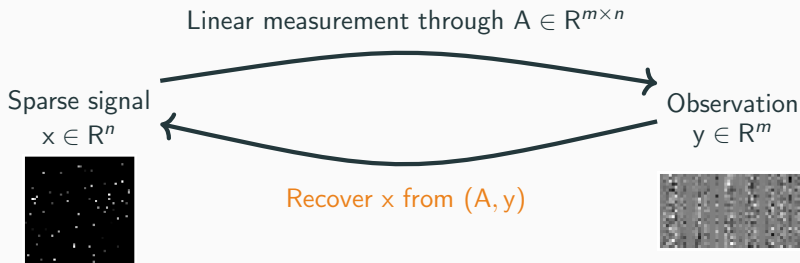
Compressive sensing



Find x such that $y \simeq Ax$

$m \ll n$: no unique solution

Compressive sensing



Find x **sparse** such that $y \simeq Ax$

$m \ll n$: no unique solution

High-dimensional statistics

Sparse GLM

Features $A \in \mathbb{R}^{m \times n}$

Targets $y \in \mathbb{R}^m$



$$\max_x \mathcal{L}(Ax, y)$$

No unique solution when $m \ll n$



$$\max_x \mathcal{L}(Ax, y) \text{ with } x \text{ sparse}$$

High-dimensional statistics

Sparse GLM

Features $A \in \mathbb{R}^{m \times n}$

Targets $y \in \mathbb{R}^m$



$$\max_x \mathcal{L}(Ax, y)$$

No unique solution when $m \ll n$



$\max_x \mathcal{L}(Ax, y)$ with x sparse

Sparse PCA

Features $A \in \mathbb{R}^{m \times n}$

Covariance $\Sigma = A^T A$



$$\max_{\|x\|_2=1} x^T \Sigma x$$

Not relevant when $m \ll n$



$\max_{\|x\|_2=1} x^T \Sigma x$ with x sparse

Machine learning

Heart disease dataset (LIBSVM)

Age	Sex	Cholesterol	Blood pressure	...	Disease
31	M	50.3 mg/dl	95 mm/hg	...	No
35	F	54.9 mg/dl	98 mm/hg	...	Yes
42	F	49.8 mg/dl	92 mm/hg	...	Yes
37	M	59.1 mg/dl	89 mm/hg	...	No
...

Machine learning

Heart disease dataset (LIBSVM)

Age	Sex	Cholesterol	Blood pressure	...	Disease
31	M	50.3 mg/dl	95 mm/hg	...	No
35	F	54.9 mg/dl	98 mm/hg	...	Yes
42	F	49.8 mg/dl	92 mm/hg	...	Yes
37	M	59.1 mg/dl	89 mm/hg	...	No
...

Data



Logistic regression

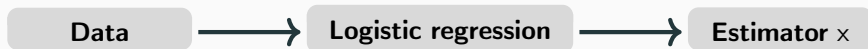


Estimator \times

Machine learning

Heart disease dataset (LIBSVM)

Age	Sex	Cholesterol	Blood pressure	...	Disease
31	M	50.3 mg/dl	95 mm/hg	...	No
35	F	54.9 mg/dl	98 mm/hg	...	Yes
42	F	49.8 mg/dl	92 mm/hg	...	Yes
37	M	59.1 mg/dl	89 mm/hg	...	No
...

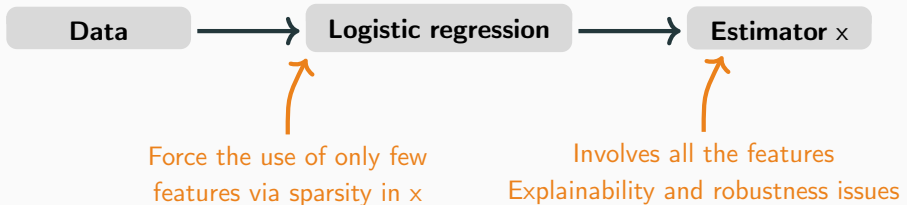


↑
Involves all the features
Explainability and robustness issues

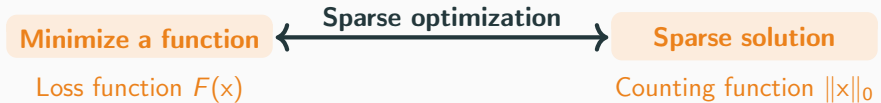
Machine learning

Heart disease dataset (LIBSVM)

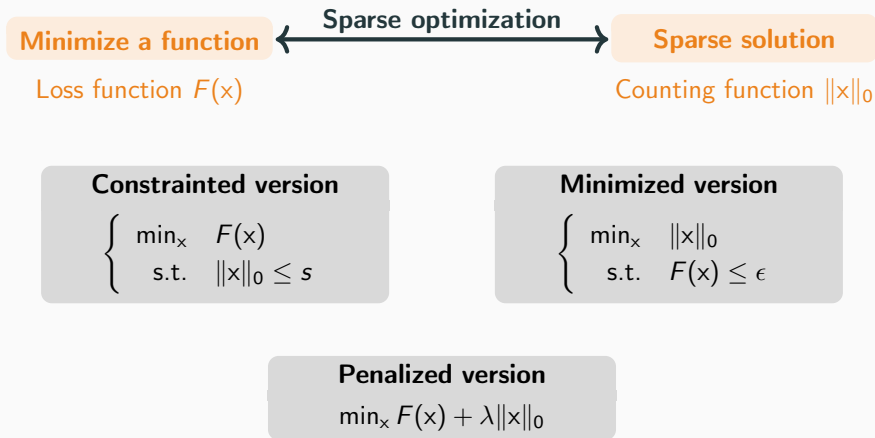
Age	Sex	Cholesterol	Blood pressure	...	Disease
31	M	50.3 mg/dl	95 mm/hg	...	No
35	F	54.9 mg/dl	98 mm/hg	...	Yes
42	F	49.8 mg/dl	92 mm/hg	...	Yes
37	M	59.1 mg/dl	89 mm/hg	...	No
...



Objective, constraint or both ?



Objective, constraint or both ?



A bit of history

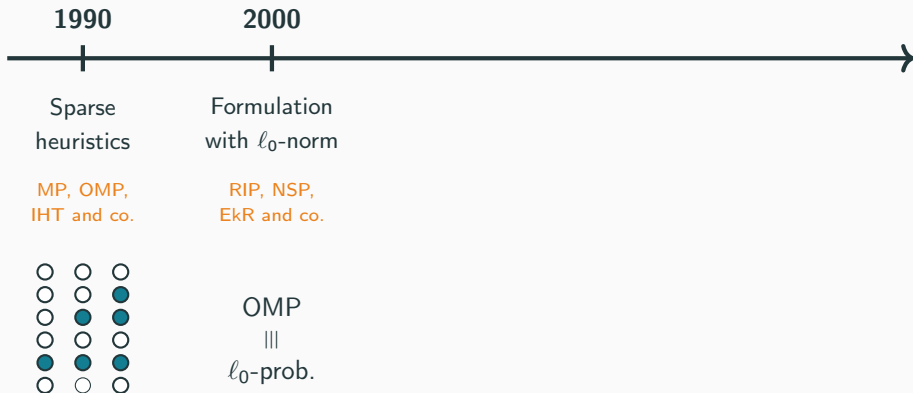
1990

Sparse
heuristics

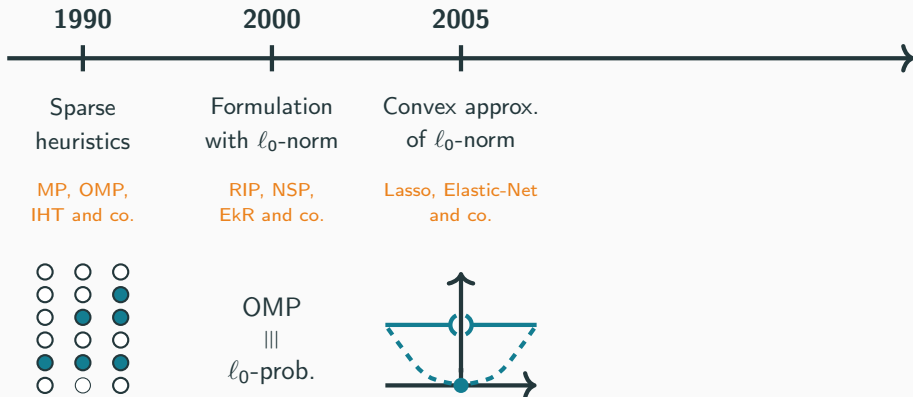
MP, OMP,
IHT and co.



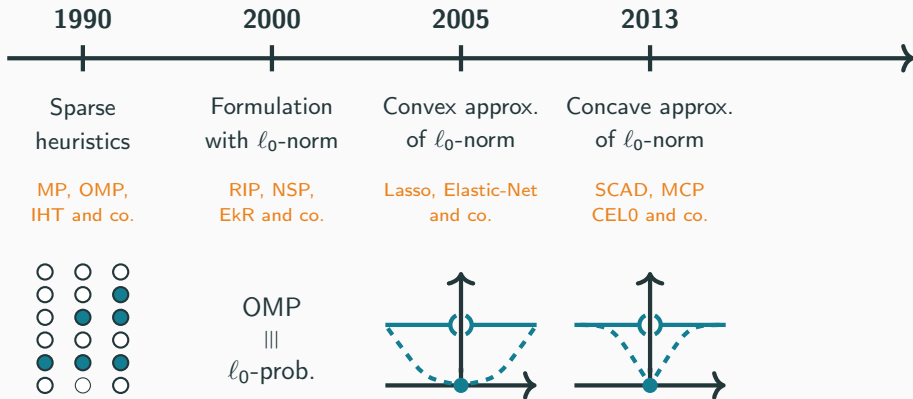
A bit of history



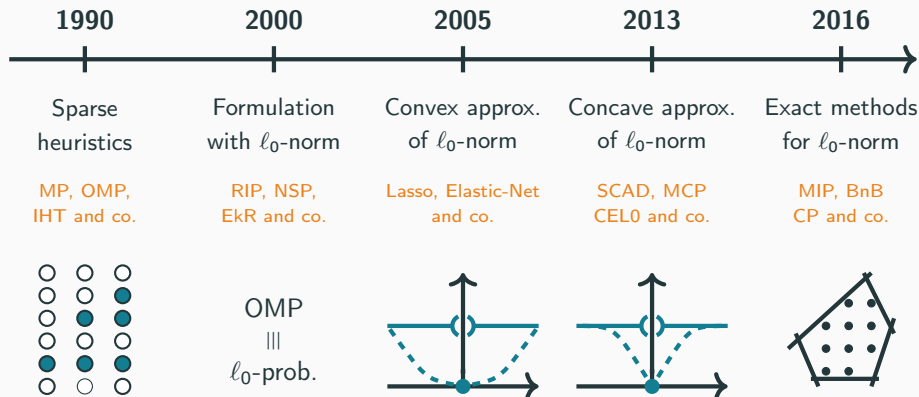
A bit of history



A bit of history

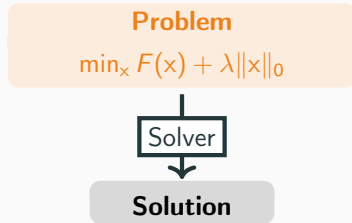


A bit of history

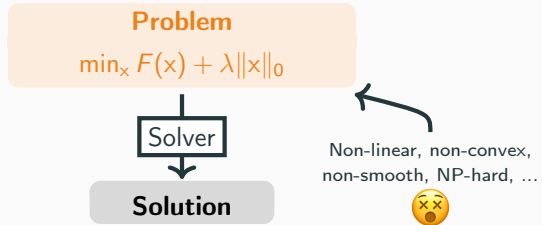


Mixed-Integer Optimization

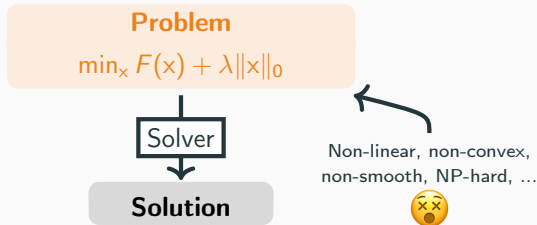
Handling the L0-norm with MIO tools



Handling the L0-norm with MIO tools



Handling the L0-norm with MIO tools



The ℓ_0 -norm **counts** the number of non-zeros in vector x

It sums the entries of the **binary** vector z satisfying some logical relation with x

We have tools to deal with such binary vectors in **MIO** !

Fitting the MIO formalism

Linearizing the ℓ_0 -norm

Real vector $x \in \mathbb{R}^n$ and binary vector $z \in \mathbb{B}^n$:

$$\|x\|_0 = 1^T z \quad \text{if} \quad x \odot (1 - z) = 0$$

Fitting the MIO formalism

Linearizing the ℓ_0 -norm

Real vector $x \in \mathbb{R}^n$ and binary vector $z \in \mathbb{B}^n$:

$$\|x\|_0 = 1^T z \quad \text{if} \quad x \odot (1 - z) = 0$$

$$\min_x F(x) + \lambda \|x\|_0$$

Fitting the MIO formalism

Linearizing the ℓ_0 -norm

Real vector $\mathbf{x} \in \mathbb{R}^n$ and binary vector $\mathbf{z} \in \mathbb{B}^n$:

$$\|\mathbf{x}\|_0 = \mathbf{1}^T \mathbf{z} \quad \text{if} \quad \mathbf{x} \odot (\mathbf{1} - \mathbf{z}) = \mathbf{0}$$

$$\min_{\mathbf{x}} F(\mathbf{x}) + \lambda \|\mathbf{x}\|_0$$



$$\min_{\mathbf{x}, \mathbf{z}} F(\mathbf{x}) + \lambda \mathbf{1}^T \mathbf{z} + H(\mathbf{x}, \mathbf{z})$$

Fitting the MIO formalism

Linearizing the ℓ_0 -norm

Real vector $x \in \mathbb{R}^n$ and binary vector $z \in \mathbb{B}^n$:

$$\|x\|_0 = 1^T z \quad \text{if} \quad x \odot (1 - z) = 0$$

$$\min_x F(x) + \lambda \|x\|_0$$



$$\min_{x,z} F(x) + \lambda 1^T z + H(x, z)$$



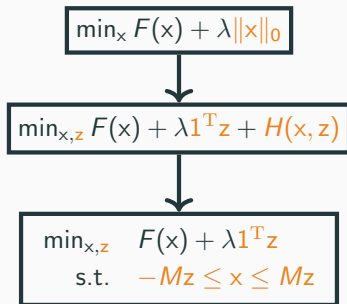
$$\begin{array}{ll} \min_{x,z} & F(x) + \lambda 1^T z \\ \text{s.t.} & -Mz \leq x \leq Mz \end{array}$$

Fitting the MIO formalism

Linearizing the ℓ_0 -norm

Real vector $x \in \mathbb{R}^n$ and binary vector $z \in \mathbb{B}^n$:

$$\|x\|_0 = 1^T z \quad \text{if} \quad x \odot (1 - z) = 0$$



Generic MIO solvers (Cplex, Gurobi, ...)

✗ Slow ✓ Generic w.r.t F/H

Specialized solvers (BnB, CP, ...)

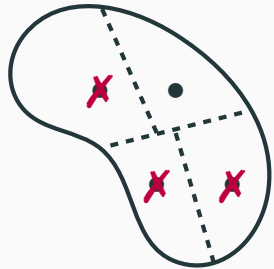
✓ Fast ✗ Restricted to some F/H

Specialized Solution Methods

Branch-and-Bound algorithms

Branch-and-Bound

“Enumerate all candidate solutions and discard sub-optimal ones.”



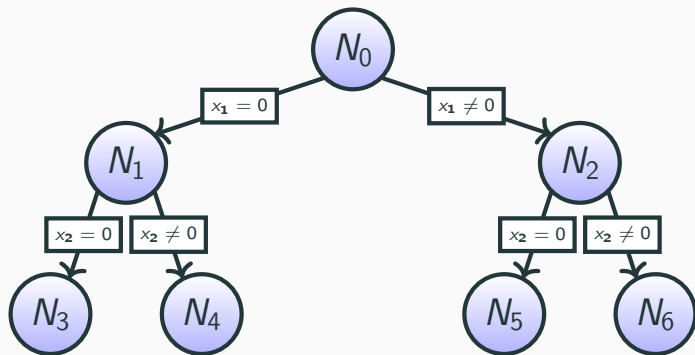
Main principles

Branching: Divide the search space

Bounding: Test whether a region can contain optimal solutions

Pruning: Discard regions without optimal solutions

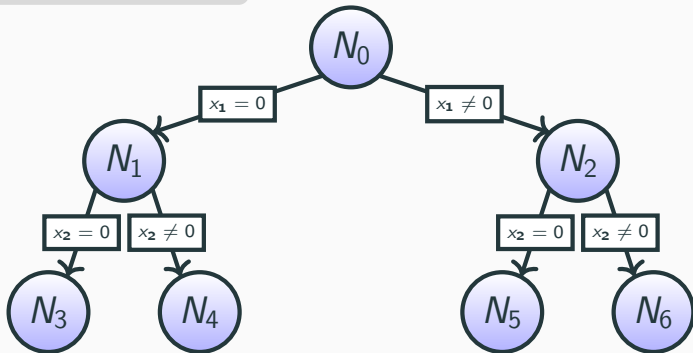
Tree exploration



Tree exploration

Observation

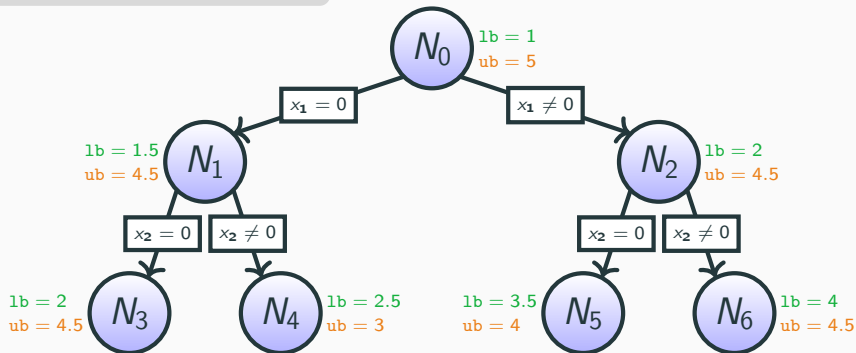
If **support of x is fixed**, then
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$
is easy to solve.



Tree exploration

Observation

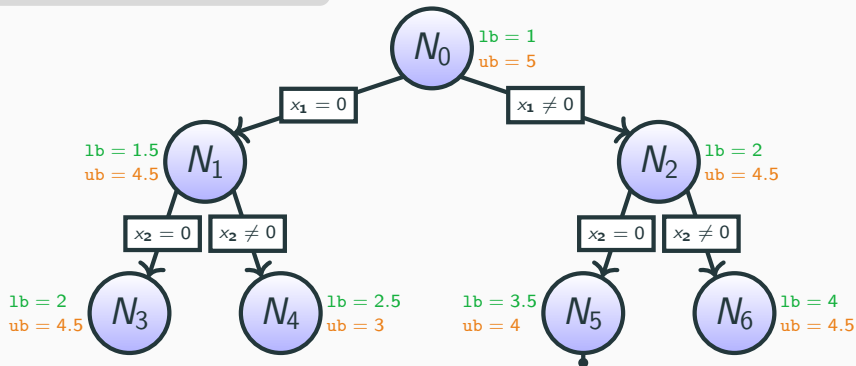
If **support of x is fixed**, then
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$
is easy to solve.



Tree exploration

Observation

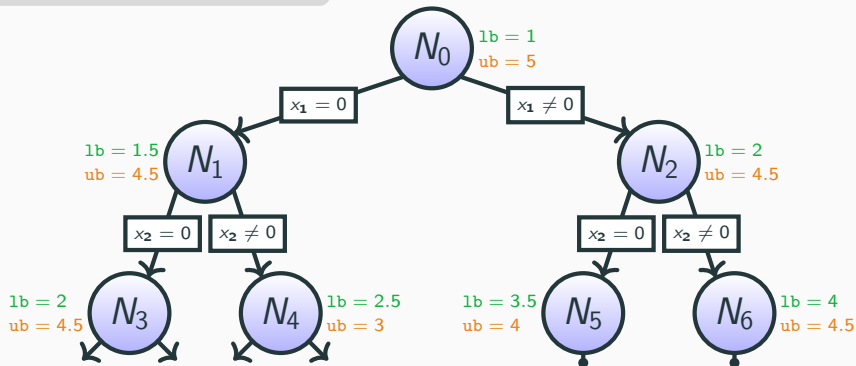
If **support of x is fixed**, then
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$
is easy to solve.



Tree exploration

Observation

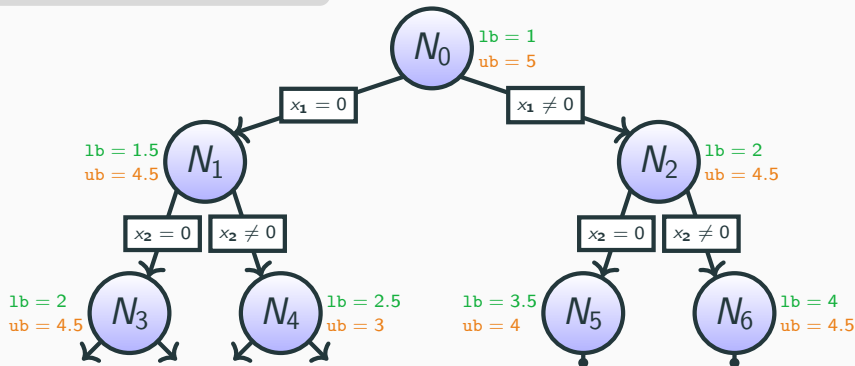
If **support of x is fixed**, then
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$
is easy to solve.



Tree exploration

Observation

If **support of x is fixed**, then
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$
is easy to solve.



All nodes explored or pruned \longrightarrow Problem solved

Node processing



Node problem

The problem at node $\nu = (S_0, S_1)$ where S_0 and S_1 are the indices of x fixed to **zero** and **non-zero** reads

$$p^\nu = \begin{cases} \min_x & F(x) + \lambda \|x\|_0 + H(x) \\ \text{s.t.} & x_{S_0} = 0, x_{S_1} \neq 0 \end{cases}$$

Node processing



Node problem

The problem at node $\nu = (S_0, S_1)$ where S_0 and S_1 are the indices of x fixed to **zero** and **non-zero** reads

$$p^\nu = \begin{cases} \min_x & F(x) + \lambda \|x\|_0 + H(x) \\ \text{s.t.} & x_{S_0} = 0, x_{S_1} \neq 0 \end{cases}$$

Task: Find lower and upper bounds on p^ν that are **tight** and **tractable to compute**

Node processing



Node problem

The problem at node $\nu = (S_0, S_1)$ where S_0 and S_1 are the indices of x fixed to **zero** and **non-zero** reads

$$p^\nu = \begin{cases} \min_x & F(x) + \lambda \|x\|_0 + H(x) \\ \text{s.t.} & x_{S_0} = 0, x_{S_1} \neq 0 \end{cases}$$

Task: Find lower and upper bounds on p^ν that are **tight** and **tractable to compute**

Upper bounding

- We just need a **feasible** solution
- Fix entries of x that are still free to zero
- Optimize the resulting problem

Node processing



Node problem

The problem at node $\nu = (S_0, S_1)$ where S_0 and S_1 are the indices of x fixed to **zero** and **non-zero** reads

$$p^\nu = \begin{cases} \min_x & F(x) + \lambda \|x\|_0 + H(x) \\ \text{s.t.} & x_{S_0} = 0, x_{S_1} \neq 0 \end{cases}$$

Task: Find lower and upper bounds on p^ν that are **tight** and **tractable to compute**

Upper bounding

- We just need a **feasible** solution
- Fix entries of x that are still free to zero
- Optimize the resulting problem

Easy to solve 🧐

Upper-bounding problem


$$\min_x F(x_{S_1}) + \lambda |S_1| + H(x_{S_1})$$

Lower bounding

Idea: Convexify a **part** of the objective function

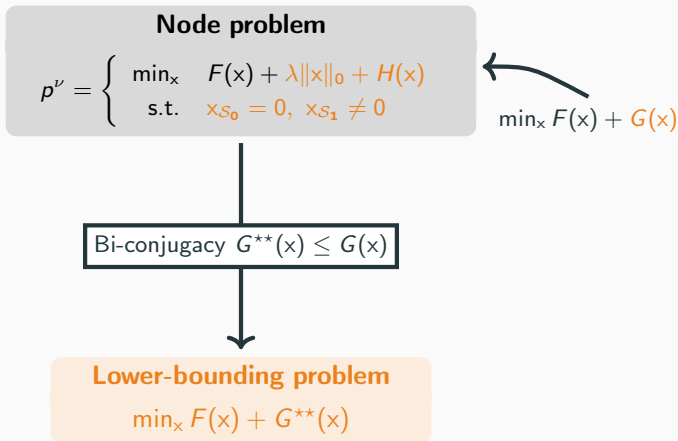
Node problem

$$p^v = \begin{cases} \min_x & F(x) + \lambda \|x\|_0 + H(x) \\ \text{s.t.} & x_{S_0} = 0, x_{S_1} \neq 0 \end{cases}$$


$$\min_x F(x) + G(x)$$

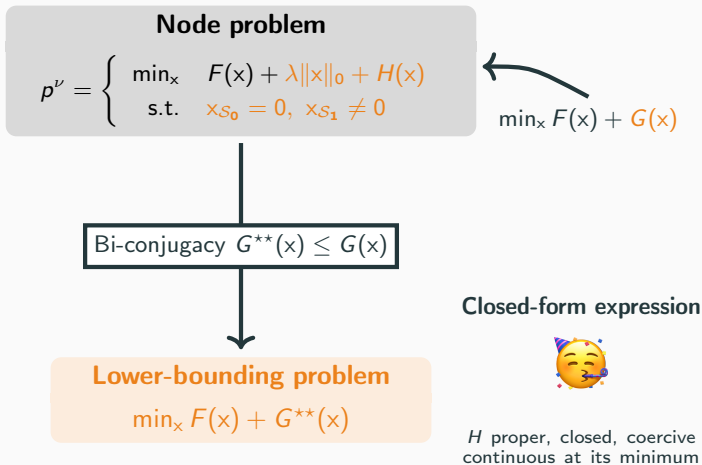
Lower bounding

Idea: Convexify a **part** of the objective function

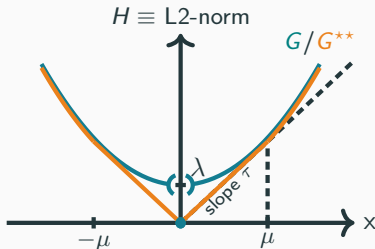
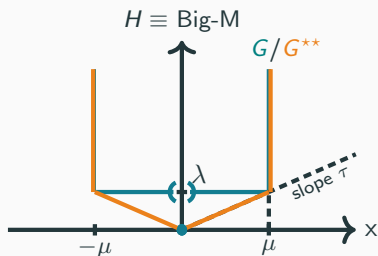


Lower bounding

Idea: Convexify a **part** of the objective function



Graphical intuition



Bi-conjugate closed-form

$$G^{**}(x) = \begin{cases} \tau|x| & \text{if } |x| \leq \mu \\ G(x) & \text{otherwise} \end{cases}$$

Let's sum up !

ℓ_0 -penalized problem

$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$

- ▶ MIP formalism
 - Linearize the ℓ_0 -norm with a binary variable
 - Big-M strategy
- ▶ Generic solvers
 - Easy solution to implement
 - Unable to exploit sparsity
 - Numerically inefficient
- ▶ Specialized Branch-and-Bound
 - Tree exploration
 - Branch by fixing support of x
 - Compute upper and lower bounds at each node
 - Leverage bi-conjugacy to compute lower bounds

Overview of Numerical Performances

Overview of numerical performances

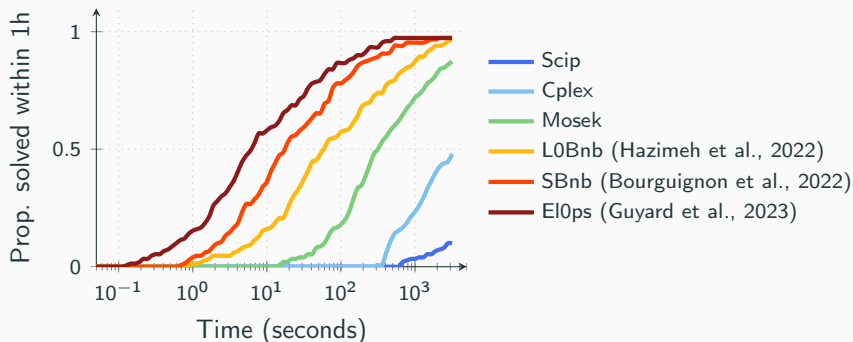
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$

Dataset : Sparse regression

F(·) : Least-squares loss

H(·) : Big-M constraints

λ : Set statistically



Overview of numerical performances

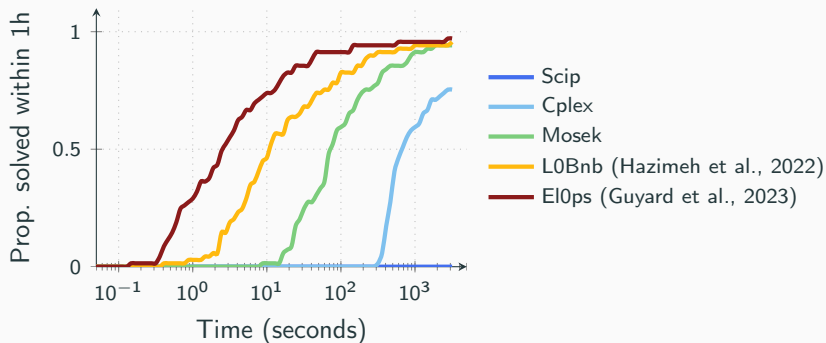
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$

Dataset : Sparse regression

F(·) : Least-squares loss

H(·) : L2-norm

λ : Set statistically



Overview of numerical performances

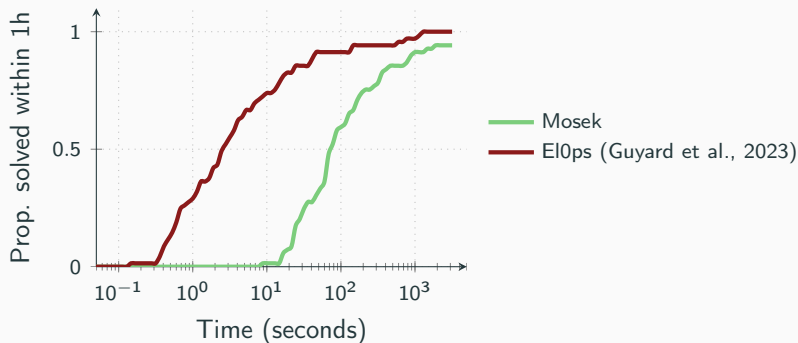
$$\min_x F(x) + \lambda \|x\|_0 + H(x)$$

Dataset : Sparse classification

F(·) : Logistic loss

H(·) : L1-norm

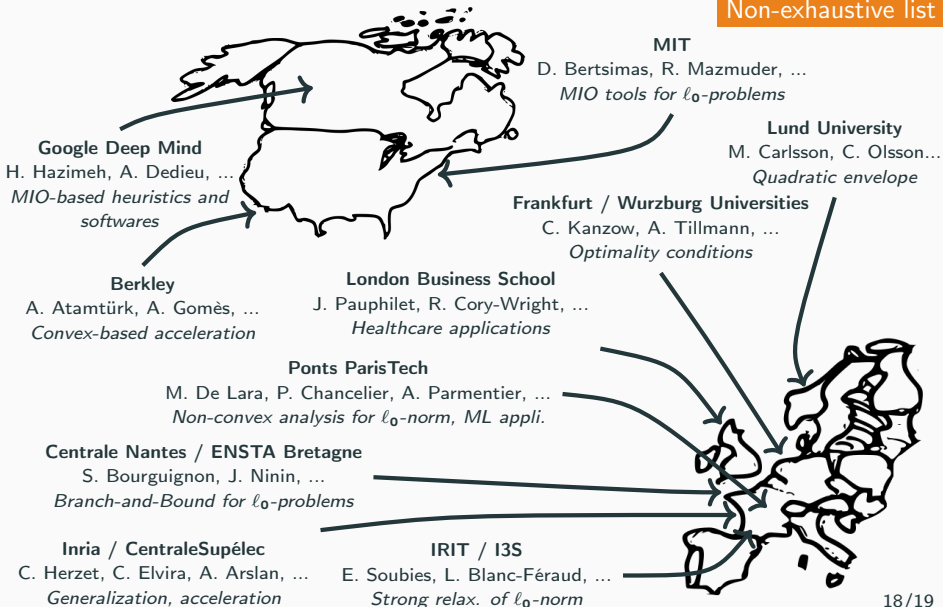
λ : Set statistically



Ongoing Research Directions

Contributors and research works

Non-exhaustive list



Take-home message

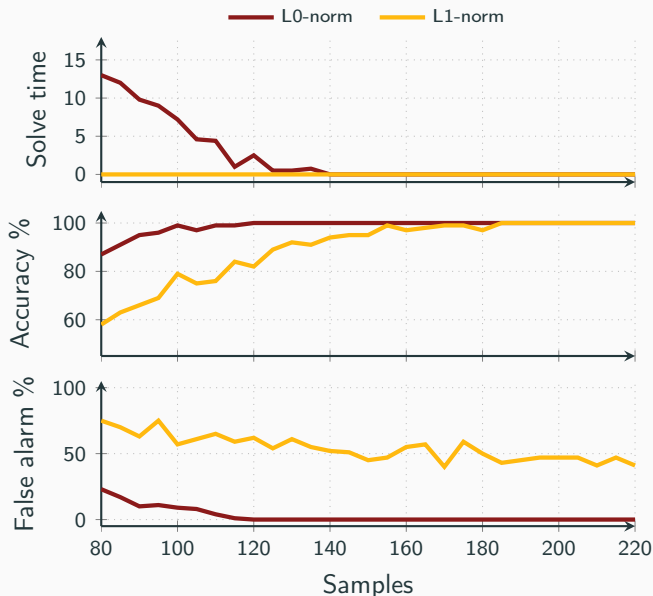
- In **some** cases, solving ℓ_0 -norm problems **exactly** worths-it
- There exists **Mixed-Integer Optimization** tools to do so
- **Structure exploitation** is the key to achieve competitive performances
- Active research area
 - Theoretical results
 - Efficiency, flexibility and accessibility of solution methods
 - Software development
 - Diffusion to other communities

Question time



Supplementary Slides

Why solving L0 problems ?



Sparse regression

$$y = Ax^\dagger + \epsilon$$

2.000 features

10 non-zeros in x^\dagger

20dB noise